

# Project 5: IMDB Movie Analysis

## Project Description:

The project aims at deriving insights from an IMDB movies dataset by answering various questions and performing analysis on the dataset. The project aims at teaching how to clean a dataset and handle the data using excel/python/R and then how to use these technologies for performing analysis and form a problem statement using 5 whys approach and then finding a solution.

## Approach:

The approach to do this project is pretty straight forward. I started out with the data cleaning process since the provided dataset is huge and such huge datasets tend to have problems like missing data, duplicate data and wrong format. After all the cleaning process, I took a look at the questions asked and found a way to answer them after performing analysis on the new clean data. For all these tasks I used Jupyter Notebook i.e. python programming language.

## Tech stack used:

- Jupyter Notebook
- Excel
- Google drive

## Insights:

- A. **Cleaning the data::** PThis is one of the most important step to perform before moving forward with the analysis. Use your knowledge learned till now to do this. (Dropping columns, removing null values, etc.)

**Your task:** Clean the data

### Steps for doing this:

- 1) Imported the libraries NumPy and pandas and imported the dataset. Took a look into the information about the dataset.
- 2) Finding percentages of null values present in the columns of the dataset. Will use this in upcoming steps.
- 3) Dropping unnecessary columns. Many of the columns were not needed so removed those.
- 4) Drop unnecessary rows using columns with high null value percentage. The gross and budget had great percentages of null values in them.

5) Removed duplicate rows.

Here are the results: size of dataset

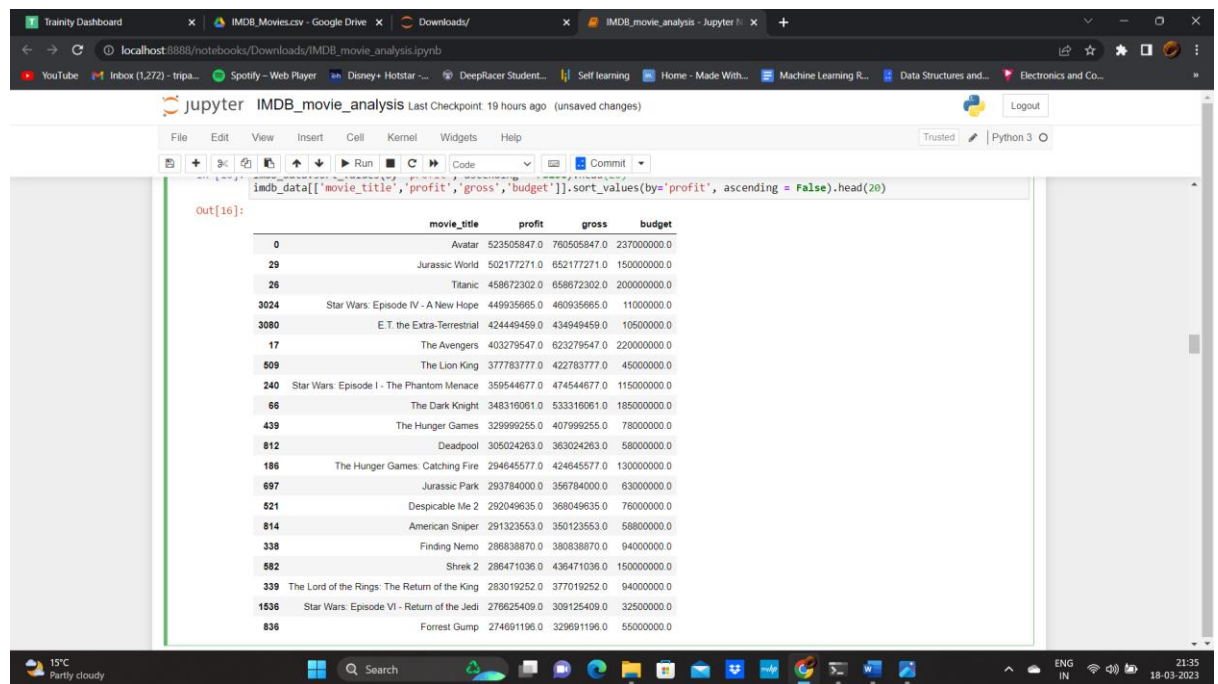
Before cleaning: 5043 x 28

After cleaning: 3849 x 13

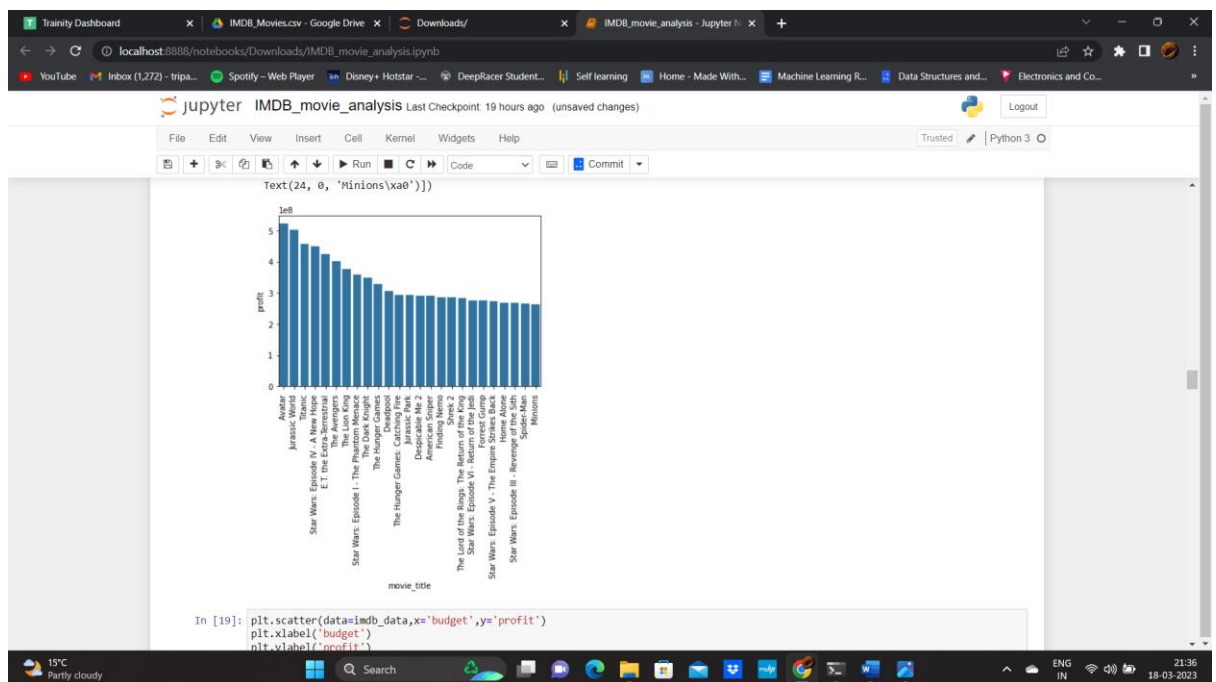
- B. **Movies with highest profit:** Create a new column called profit which contains the difference of the two columns: gross and budget. Sort the column using the profit column as reference. Plot profit (y-axis) vs budget (x-axis) and observe the outliers using the appropriate chart type.

**Your task:** Find the movies with the highest profit?

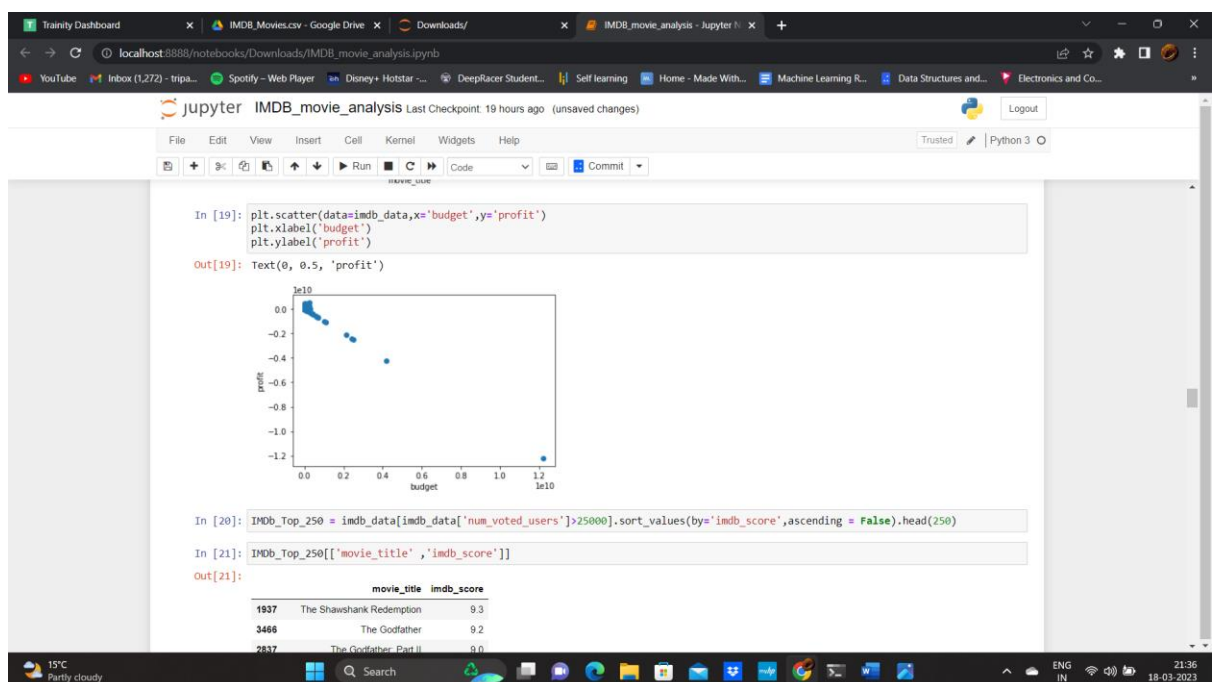
**The following are the top 20 movies with highest profit:**



Bar graph representation of top 20 movies:



The outliers can be observed using scatter plot:



- C. **Top 250:** Create a new column IMDb\_Top\_250 and store the top 250 movies with the highest IMDb Rating (corresponding to the column: imdb\_score). Also make sure that for all of these movies, the num\_voted\_users is greater than 25,000. Also add a Rank column containing the values 1 to 250 indicating the ranks of the

corresponding films.

Extract all the movies in the IMDB\_Top\_250 column which are not in the English language and store them in a new column named Top\_Foreign\_Lang\_Film. You can use your own imagination also!

**Your task:** Find IMDB Top 250

The screenshot shows a Jupyter Notebook interface with two code cells. The first cell filters the IMDB\_Top\_250 dataset to show only non-English movies. The second cell filters the dataset to show only English movies. The output of the first cell is a table with 250 rows and 4 columns: movie\_title, IMDB\_Top\_250, num\_voted\_users, and Rank. The output of the second cell is a table with 250 rows and 7 columns: director\_name, num\_critic\_for\_reviews, director\_facebook\_likes, gross, genres, actor\_1\_name, and movie\_title.

```
In [24]: IMDB_Top_250[IMDB_Top_250['language'] != 'English']
IMDB_Top_250[['movie_title', 'IMDB_Top_250', 'num_voted_users', 'Rank']]

Out[24]:
```

	movie_title	IMDB_Top_250	num_voted_users	Rank
1937	The Shawshank Redemption	9.3	1689764	1.0
3466	The Godfather	9.2	1155770	2.0
2837	The Godfather: Part II	9.0	790926	3.0
66	The Dark Knight	9.0	1676169	4.0
4498	The Good, the Bad and the Ugly	8.9	503509	5.0
...	...	...	...	...
4931	Once	7.9	90827	246.0
2605	Crouching Tiger, Hidden Dragon	7.9	217740	247.0
3029	The Fighter	7.9	275869	248.0
2177	Edward Scissorhands	7.9	357581	249.0
2487	My Fair Lady	7.9	66959	250.0

```
In [25]: Top_Foreign_Lang_Film = IMDB_Top_250[IMDB_Top_250['language'] != 'English']
Top_Foreign_Lang_Film

Out[25]:
```

	director_name	num_critic_for_reviews	director_facebook_likes	gross	genres	actor_1_name	movie_title
4498	Sergio Leone	181.0	0.0	61000000.0	Western	Clint Eastwood	The Good, the Bad and the Ugly

The following picture shows movie names and IMDB\_Top\_250 column with top imdb scores along with rank column from 1-250

Top\_foreign\_language\_Films:

Rank	Top_Foreign_Lang_Film	language
4488	5.0 The Good, the Bad and the Ugly	Italian
4747	17.0 Seven Samurai	Japanese
4029	20.0 City of God	Portuguese
2373	23.0 Spirited Away	Japanese
4259	35.0 The Lives of Others	German
4921	39.0 Children of Heaven	Persian
2323	47.0 Princess Mononoke	Japanese
2970	49.0 Das Boot	German
4105	57.0 Oldboy	Korean
4689	58.0 A Separation	Persian
1329	59.0 Baahubali: The Beginning	Telugu
1298	61.0 Amélie	French
2734	65.0 Metropolis	German
4033	73.0 The Hunt	Danish
2829	81.0 Downfall	German
2551	86.0 Pan's Labyrinth	Spanish
4000	97.0 The Secret in Their Eyes	Spanish
3550	98.0 Incendies	French
2047	100.0 How's Moving Castle	Japanese
2830	110.0 The Sea Inside	Spanish
2914	119.0 Tae Guk Gi: The Brotherhood of War	Korean
4451	127.0 The Celebration	Danish
3553	143.0 Elite Squad	Portuguese
3423	147.0 Akira	Japanese
4267	152.0 Amores Perros	Spanish

D. Best Directors: TGroup the column using the director\_name column.

Find out the top 10 directors for whom the mean of imdb\_score is the highest and store them in a new column top10director. In case of a tie in IMDb score between two directors, sort them alphabetically.

Your task: Find the best directors

```

In [148]: top10director = movies.groupby('director_name').imdb_score.mean().sort_values(ascending = False).head(10)
          movies['imdb_score_mean'] = movies['imdb_score'].mean()
          top10director

Out[148]: director_name
Charles Chaplin      8.600000
Tony Kaye             8.600000
Alfred Hitchcock     8.500000
Ron Fricke           8.500000
Damien Chazelle      8.500000
Majid Majidi         8.500000
Sergio Leone         8.433333
Christopher Nolan    8.425000
S.S. Rajamouli       8.400000
Marius A. Markevicius 8.400000
Name: imdb_score, dtype: float64

```

The above screenshot shows list of top 10 directors based on the mean of imdb scores.

E. Popular Genres: Perform this step using the knowledge gained while performing previous steps.

Your task: Find popular genres

```

Majid Majidi      8.500000
Sergio Leone      8.433333
Christopher Nolan  8.425000
Aghar Farhadi     8.400000
Richard Marquand  8.400000
Name: imdb_score, dtype: float64

In [28]: top10genres = imdb_data.groupby('genres').imdb_score.mean().sort_values(ascending = False).head(10)
top10genres
Out[28]:
genres
Crime|Drama|Fantasy|Mystery      8.50
Adventure|Animation|Drama|Family|Musical  8.50
Action|Adventure|Drama|Fantasy|War      8.40
Adventure|Animation|Fantasy      8.40
Adventure|Drama|Thriller|War      8.40
Documentary|War      8.30
Biography|Drama|History|Music      8.30
Adventure|Animation|Comedy|Drama|Family|Fantasy  8.30
Documentary|Drama|Sport      8.30
Adventure|Drama|War      8.25
Name: imdb_score, dtype: float64

In [29]: Meryl_Streep = imdb_data[imdb_data['actor_1_name']=='Meryl Streep']
Leo_Caprio = imdb_data[imdb_data['actor_1_name']=='Leonardo DiCaprio']
Brad_Pitt = imdb_data[imdb_data['actor_1_name']=='Brad Pitt']

In [30]: Combined = Meryl_Streep.append([Leo_Caprio,Brad_Pitt])
Combined
792 Patrick Gilmore      98.0      0.0  26288320.0  Adventure|Animation|Comedy|Drama|Family|Fantas...  Brad Pitt

```

Following are the top 10a genres .

- F. Charts: Create three new columns namely, Meryl\_Streep, Leo\_Caprio, and Brad\_Pitt which contain the movies in which the actors: 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt' are the lead actors. Use only the actor\_1\_name column for extraction. Also, make sure that you use the names 'Meryl Streep', 'Leonardo DiCaprio', and 'Brad Pitt' for the said extraction.

Append the rows of all these columns and store them in a new column named Combined.

Group the combined column using the actor\_1\_name column.

Find the mean of the num\_critic\_for\_reviews and num\_users\_for\_review and identify the actors which have the highest mean.

Observe the change in number of voted users over decades using a bar chart. Create a column called decade which represents the decade to which every movie belongs to. For example, the title\_year year 1923, 1925 should be stored as 1920s. **Sort the column based on the column decade, group it by decade and find the sum of users voted in each decade. Store this in a new data frame called df\_by\_decade.**

Your task: Find the critic-favorite and audience-favorite actors

1490	Terrence Mallick	584.0	0.0	13303319.0	Drama Fantasy	E
1722	Andrew Dominik	273.0	181.0	3904982.0	Biography Crime Drama History Western	E
2204	Alejandro G. Iñárritu	285.0	0.0	34300771.0	Drama	E
2333	Angelina Jolie Pitt	131.0	11000.0	531009.0	Drama Romance	E

```

In [31]: Combined.groupby('actor_1_name').num_critic_for_reviews.mean()

Out[31]: actor_1_name
Brad Pitt          245.000000
Leonardo DiCaprio 330.190476
Meryl Streep       181.454545
Name: num_critic_for_reviews, dtype: float64

In [32]: Combined.num_user_for_reviews=Combined.num_user_for_reviews.astype('int')

In [33]: Combined.groupby('actor_1_name').num_user_for_reviews.mean()

Out[33]: actor_1_name

```

```

In [31]: Combined.groupby('actor_1_name').num_critic_for_reviews.mean()

Out[31]: actor_1_name
Brad Pitt          245.000000
Leonardo DiCaprio 330.190476
Meryl Streep       181.454545
Name: num_critic_for_reviews, dtype: float64

In [32]: Combined.num_user_for_reviews=Combined.num_user_for_reviews.astype('int')

In [33]: Combined.groupby('actor_1_name').num_user_for_reviews.mean()

Out[33]: actor_1_name
Brad Pitt          742.352941
Leonardo DiCaprio 914.476190
Meryl Streep       297.181818
Name: num_user_for_reviews, dtype: float64

In [34]: min_year = imdb_data.title_year.min()

In [35]: max_year=imdb_data['title_year'].max()
int(max_year)+(10 -int(max_year)%10)

Out[35]: 2020

```

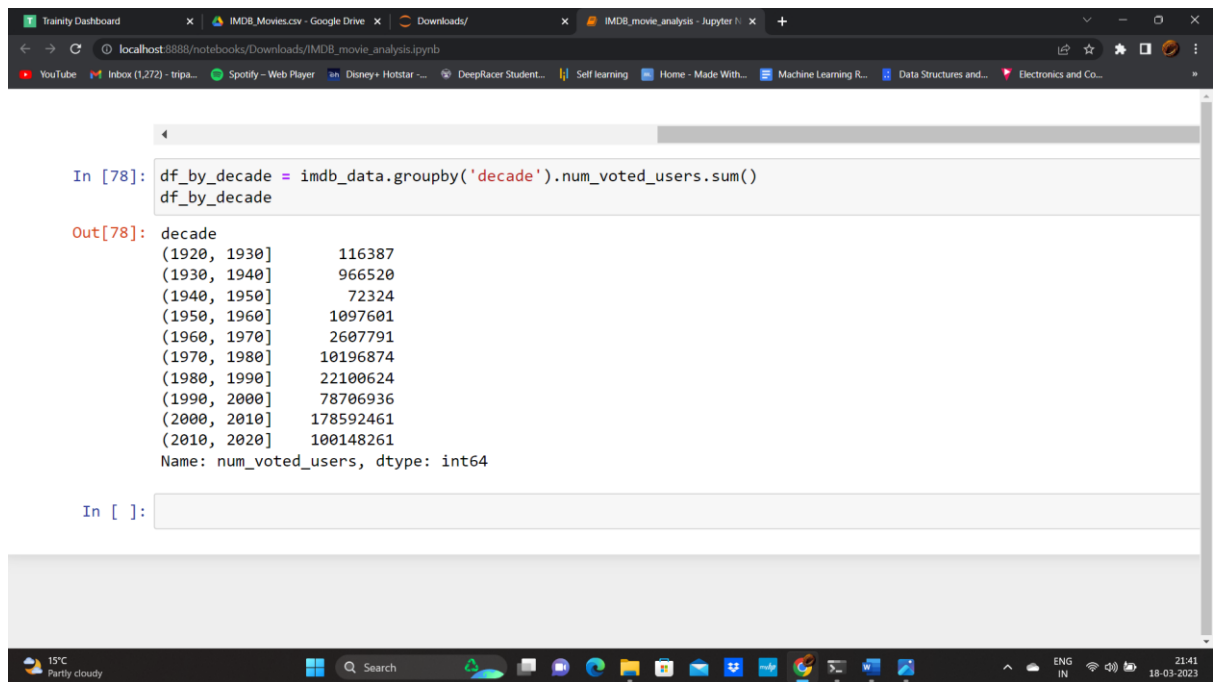
Shows the

- 1) Critic favourite: Leonardo DiCaprio
- 2) Audience favourite : Leonardo DiCaprio



Observe the change in number of voted users over decades using a bar chart. Create a column called decade which represents the decade to which every movie belongs to. For example, the title\_year year 1923, 1925 should be stored as 1920s. Sort the column based on the column decade, group it by decade and find the sum of users voted in each decade. Store this in a new data frame called df\_by\_decade.

The following ss will show the sum of user votes as per decade:



The screenshot shows a Jupyter Notebook interface with the following content:

```
In [78]: df_by_decade = imdb_data.groupby('decade').num_voted_users.sum()
df_by_decade
```

```
Out[78]: decade
(1920, 1930]    116387
(1930, 1940]    966520
(1940, 1950]     72324
(1950, 1960]   1097601
(1960, 1970]   2607791
(1970, 1980]  10196874
(1980, 1990]  22100624
(1990, 2000]  78706936
(2000, 2010] 178592461
(2010, 2020] 100148261
Name: num_voted_users, dtype: int64
```

```
In [ ]:
```

The notebook is running in a browser at localhost:8888. The taskbar at the bottom shows the system time as 21:41 on 18-03-2023.

### Result:

Doing this project involved many challenges which made me understand everything and work on it. Using python and its popular libraries for data cleaning and analysis is what in which now I am confident.



